



## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

<b>(51) International Patent Classification <sup>7</sup> :</b> <b>G06F 9/00</b>	<b>A2</b>	<b>(11) International Publication Number:</b> <b>WO 00/62157</b> <b>(43) International Publication Date:</b> 19 October 2000 (19.10.00)
<b>(21) International Application Number:</b> PCT/EP00/03204 <b>(22) International Filing Date:</b> 11 April 2000 (11.04.00)  <b>(30) Priority Data:</b> 60/129,301 14 April 1999 (14.04.99) US 09/481,771 11 January 2000 (11.01.00) US  <b>(71) Applicant:</b> KONINKLIJKE PHILIPS ELECTRONICS N.V. [NL/NL]; Groenewoudseweg 1, NL-5621 BA Eindhoven (NL).  <b>(72) Inventor:</b> ISHAM, Karl, M.; Prof. Holstlaan 6, NL-5656 AA Eindhoven (NL).  <b>(74) Agent:</b> GRAVENDEEL, Cornelis; Internationaal Octrooi- bureau B.V., Prof. Holstlaan 6, NL-5656 AA Eindhoven (NL).		<b>(81) Designated States:</b> JP, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).  <b>Published</b> <i>Without international search report and to be republished upon receipt of that report.</i>
<b>(54) Title:</b> METHOD FOR DYNAMIC LOANING IN RATE MONOTONIC REAL-TIME SYSTEMS  <b>(57) Abstract</b>  A method and apparatus are disclosed for sharing execution capacity among tasks executing in a real-time computing system. The present invention extends RMA techniques for characterizing system timing behavior and designing real-time systems. A high priority task having hard deadlines is paired with a lower priority task having soft deadlines. During an overload condition, the higher priority task can dynamically borrow execution time from the execution capacity of the lower priority task without affecting the schedulability of the rest of the system. The higher priority task is bolstered in a proportion to the capacity borrowed from the lower priority task, so that the combined utilization of the two tasks remains constant. The period of the degraded task is increased to compensate for the execution time that was loaned to the higher priority task. In addition, the priority of the lower priority task is modified to match the new period.		

BEST AVAILABLE COPY

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav	TM	Turkmenistan
BF	Burkina Faso	GR	Greece		Republic of Macedonia	TR	Turkey
BG	Bulgaria	HU	Hungary	ML	Mali	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MN	Mongolia	UA	Ukraine
BR	Brazil	IL	Israel	MR	Mauritania	UG	Uganda
BY	Belarus	IS	Iceland	MW	Malawi	US	United States of America
CA	Canada	IT	Italy	MX	Mexico	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NE	Niger	VN	Viet Nam
CG	Congo	KE	Kenya	NL	Netherlands	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NO	Norway	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's	NZ	New Zealand		
CM	Cameroon		Republic of Korea	PL	Poland		
CN	China	KR	Republic of Korea	PT	Portugal		
CU	Cuba	KZ	Kazakstan	RO	Romania		
CZ	Czech Republic	LC	Saint Lucia	RU	Russian Federation		
DE	Germany	LI	Liechtenstein	SD	Sudan		
DK	Denmark	LK	Sri Lanka	SE	Sweden		
EE	Estonia	LR	Liberia	SG	Singapore		

## Method for dynamic loaning in rate monotonic real-time systems.

5           The present invention relates to the timing behavior of real-time computing systems, and more particularly, to a method and apparatus for loaning execution capacity among tasks executing in a real-time computing system.

10           Real-time systems differ from other forms of computing in that they must be temporally correct as well as logically correct. Such systems are developed to satisfy the following three primary criteria: guaranteed timing deadlines, fast response times, and stability in overload. Schedulability describes the capacity of a system. A system that is schedulable can meet all of its critical timing deadlines. Latency describes the responsiveness  
15 of a system. In a real-time system, it is the worst-case system response time to events that matters. Stability in overload means the system is able to meet its critical deadlines even if all deadlines cannot be met.

          One of the most useful models available for developing real-time computing systems is Rate Monotonic Analysis (RMA). RMA provides a mathematical framework for  
20 reasoning about system timing behavior and provides an engineering basis for designing real-time systems. RMA was first disclosed in Liu & Layland, "Scheduling Algorithms for Multi-Programming in a Hard Real-Time Environment," Journal of the Ass'n of Computing Machinery (ACM) 20, 1, 40-61 (January, 1973), incorporated by reference herein. Generally, Liu and Layland demonstrated that a set of  $n$  periodic tasks with deadlines at the end of their  
25 periods will meet their deadlines if they are arranged in priority according to their periods, and they meet a schedulability bound test. Since this paper, RMA has evolved into a collection of methods that have extended the original theory from its original form. The basic RMA real-time guarantee, however, has remained unchanged. For a general discussion of these collections of methods, see, for example, Klein et al, A Practitioner's Handbook for

Real-Time Analysis: Guide to Rate Monotonic Analysis for Real-time Systems (Kluwer Academic Publishing, ISBN 0-7923-9361-9, 1993), incorporated by reference herein.

Methods currently exist to assess spare capacity and dynamically change the priority of tasks to alter the scheduling policy of the system. Spare capacity is the amount of execution time that can be added to the response of an event while preserving the schedulability of lower priority events. A related method, eliminating overrun, computes the amount of resource usage that must be eliminated to allow an event to meet its deadlines.

See, Klein et al, A Practitioner's Handbook for Real-Time Analysis: Guide to Rate Monotonic Analysis for Real-time Systems, Chapter 4, Group 3 (Kluwer Academic Publishing, ISBN 0-7923-9361-9, 1993). Changing the priority of a task is used in various synchronization protocols to avoid priority inversion when multiple tasks share common data.

The validity of Rate Monotonic Analysis depends on preparing for the worst-case in terms of individual event timing and the concurrence of events. In other words, the system will be capable of meeting all of its deadlines if all worst-case events can be handled simultaneously. While mathematically certain, assuming worst-case timings and coincidences required by the RMA analysis may result in harsh feature/performance tradeoffs for the usually limited hardware computing capacity found in a typical real-time system. While the worst-case high priority event may occur very infrequently, the execution time for the worst-case high priority event must still be accounted for in RMA calculations. The capacity allocated to handle this rare worst-case event at high priority then gets traded off against the capacity available to lower priority events.

If RMA calculations show that the system is no longer schedulable with these maximum execution requirements, the designer is left with few options. Usually, the worst-case events will have to be redesigned to reduce their execution times. If the execution times cannot be reduced, the designer is left with discounting the offending event as an overload, provided the occasional timing violations to lower priority events can be tolerated. An overload condition is an acceptable alternative only if all lower priority events have soft deadlines (can tolerate a missed deadline occasionally). Unfortunately, this is rarely the case. In a system where such an overload situation exists, the real-time criteria of stability in overload is violated.

Currently, RMA only provides the tools necessary to determine if a system is schedulable. In other words, RMA will only tell you if a given design will work or not. For example, the method of assessing spare capacity is useful for determining how much

execution time can be safely added to an event before timing requirements are broken. Likewise, RMA provides a method for determining how much execution time must be eliminated to meet timing requirements. These methods will prove useful for providing target values for redesign, but will not help if the execution times are intractable. There is no  
5 current method that allows a designer to gracefully handle intractable overloads.

Generally, a method and apparatus are disclosed for sharing execution capacity among tasks executing in a real-time computing system by pairing a high priority  
10 task having hard deadlines with a lower priority task having soft deadlines. The present invention extends RMA techniques for characterizing system timing behavior and designing real-time systems. During an overload condition, the higher priority task can dynamically borrow execution time from the execution capacity of the lower priority task without affecting the schedulability of the rest of the system. The higher priority task is bolstered in  
15 proportion to the capacity borrowed from the lower priority task, so that the combined utilization of the two tasks remains constant.

According to another aspect of the invention, the period of the degraded task is increased to compensate for the execution time that was loaned to the higher priority task. Thus, events for the degraded task may still be assigned according to the original execution  
20 budget, but are allowed a longer time to complete due to the increased work in the borrowing higher priority task. In addition, the priority of the lower priority task is modified to match the new period, in accordance with the Rate Monotonic Scheduling algorithm.

The present invention isolates effects of an overload to a particular task (i.e., the degraded task). In addition, the manner in which execution capacity is borrowed does not  
25 diminish the original execution budget of the degraded task, since the period of the degraded task is lengthened to compensate for the borrowed time. Thus, the overloaded system takes longer to perform non-critical events but the amount of work it will accept at any one time remains the same. In addition, the present invention provides better utilization of computing resources by borrowing non-overload execution time against the capacity of a rarely used  
30 task. In this case, a low priority task may be prepared to handle an intermittent event with soft deadlines. The present invention permits a higher priority task to borrow some of the capacity for its normal operation from the intermittent task.

A more complete understanding of the present invention, as well as further features and advantages of the present invention, will be obtained by reference to the following detailed description and drawings.

5

FIG. 1 illustrates a real-time computing system in accordance with the present invention; and

FIG. 2 illustrates a state diagram of the capacity loaning process of FIG. 1.

10

The present invention provides a method for pairing a high priority task having hard deadlines with a lower priority task with soft deadlines. During an overload, the higher priority task can dynamically borrow execution time from the execution capacity of the lower priority task without affecting the schedulability of the rest of the system. In this manner, the higher priority task can be bolstered in a proportional manner, so that the combined utilization of the two tasks remains constant.

15

According to another aspect of the present invention, the period of the degraded loaning task is lengthened to compensate for the borrowed time. In addition, the priority of the degraded task is modified to match the new period. Furthermore, the performance of the loaning task is degraded for the period of the borrowing task, extended over the whole period of the loaning task.

20

The present invention provides a means to gracefully degrade system performance during overload while still providing critical schedulability guarantees. For intractable overloads, this may be the only practical method of maintaining stability. Secondly, the task to be degraded can be identified during the design stage, which is an important consideration from the standpoint of predictability and isolation. Real-time problems are notorious for their ability to elude capture in the field. If an overload event coincided with a different lower priority event every time it manifested itself, it would be nearly impossible to correctly identify and correct.

25

The present invention isolates the effects of an overload to a particular task. In addition, the means used by the present invention to borrow execution capacity does not diminish the original execution budget of the degraded task, since the period of the degraded task is lengthened to compensate for the borrowed time. Thus, the overloaded system takes longer to perform non-critical events but the amount of work it will accept at any one time

30

remains the same. Finally, the method of the present invention can be used to provide better utilization of computing resources by borrowing non-overload execution time against the capacity of a rarely used task. In this case, a low priority task may be prepared to handle an intermittent event with soft deadlines. The present invention permits a higher priority task to borrow some of the capacity for its normal operation from the intermittent task.

As previously indicated, Liu and Layland demonstrated that a set of  $n$  periodic tasks with deadlines at the end of their periods will meet their deadlines if they are arranged in priority according to their periods, and they meet a schedulability bound test. For a detailed discussion of the schedulability bound test, see Liu & Layland, "Scheduling Algorithms for Multi-Programming in a Hard Real-Time Environment," Journal of the Association of Computing Machinery (ACM) 20, 1, 40-61 (January, 1973), incorporated by reference herein.

Generally, the schedulability bound test states that a set of  $n$  independent periodic tasks scheduled by the rate monotonic algorithm will always meet its deadlines, for all task phasings, if:

$$\frac{C_1}{T_1} + \dots + \frac{C_n}{T_n} \leq U(n) = n \left( 2^{\frac{1}{n}} - 1 \right)$$

where,

$C_i$  = worst-case task execution time of task  $i$ ,

$T_i$  = period of task  $i$ , and

$U(n)$  = utilization for  $n$  tasks.

As previously indicated, the present invention pairs two tasks, a critical task with hard deadlines and another task of lower priority with soft deadlines. Consider a pair of tasks whose total utilization remains constant, thereby allowing RMA to be applied to guarantee schedulability. Such a pair might exist in a real-time system in order to satisfy two different quality of service levels. The task pair relationship can be expressed as follows:

$$\frac{C_u}{T_u} + \frac{C_r}{T_r} = U \quad \text{Equation (1)}$$

where,

$C_u$  = worst-case task execution time of task  $u$ ,

$T_u$  = period of task<sub>u</sub>,

$C_r$  = worst-case task execution time of task<sub>r</sub>,

$T_r$  = period of task<sub>r</sub>, and

$U$  = utilization for both tasks.

Thus, the utilization of one task in the pair may be increased if the utilization of the other task in the pair is proportionally decreased to maintain a constant utilization,  $U$ . In this manner, execution time can be borrowed from one task to supplement the execution of the other task. The techniques of the present invention are useful, for example, where events are validated against the available execution time,  $C_u$  or  $C_r$ , before they are assigned to a servicing task, task<sub>u</sub> or task<sub>r</sub>. It follows that temporary overloads may be assigned to the higher priority task without sacrificing RMA guarantees. The utilization loan can be expressed by the equation:

$$\frac{C_u + N_u}{T_u} + \frac{C_r - N_r}{T_r} = U$$

$$\frac{C_u}{T_u} + \frac{N_u}{T_u} + \frac{C_r}{T_r} - \frac{N_r}{T_r} = U$$

$$\frac{N_u}{T_u} - \frac{N_r}{T_r} = U - \frac{C_u}{T_u} - \frac{C_r}{T_r}$$

$$\frac{N_u}{T_u} - \frac{N_r}{T_r} = 0 \quad (\text{from equation (1)})$$

$$\frac{N_u}{T_u} = \frac{N_r}{T_r}$$

$$N_u = \frac{N_r \cdot T_u}{T_r} \quad \text{Equation (2)}$$

where,

$N_r$  = amount of execution time to borrow from task<sub>r</sub>, where  $N_r < C_r$ ; and

$N_u$  = amount of execution time available to loan to task<sub>u</sub>.

The utilization terms for the two tasks in question are isolated from the utilization bound equation and equated to their combined utilization. In the above discussion, task<sub>u</sub> refers to the higher priority (urgent) task while task<sub>r</sub> refers to the lower priority (routine) task with soft deadlines. For the RMA utilization bound test to remain valid, the sum of these two tasks must remain constant during the application of the method.

Given that the utilization,  $U$ , remains constant, the utilization of the higher priority task may be increased as long as there is a proportional decrease in the utilization of the lower priority task. The proportional change guarantees the schedulability of the rest of



the system, since the net effect on the utilization bound equation is the same before and after the change. Equation (2) applies this proportional change and is solved for the amount of execution time the higher priority task may borrow.

While adding execution time to a task may have an application in handling overloads, it is usually not practical to actually remove execution time from another task to compensate for it. The execution times of most events at runtime are fixed, and therefore cannot be limited. This is solved by assuming that the original execution budget is fixed and letting the period of the task change instead. As discussed hereinafter, Equations (3) and (4) deal with this change in the loaning task's period. Since the period of a task also determines its priority in a Rate Monotonic environment, the priority of the task must change as well, using the usual set priority functions available in most real-time operating systems.

To compensate for the borrowed execution time, the period of the loaning task may be changed instead of limiting the execution time of the task. In this manner, the level of service in the loaning task is degraded gracefully. Events may still be assigned according to the original execution budget,  $C_r$ , but are allowed a longer time to complete due to the increased work in the borrowing task. The graceful degradation may be implemented as follows:

$$\begin{aligned}\frac{C_r}{T_n} &= \frac{C_r - N_r}{T_r} \\ T_n &= \frac{C_r}{\left(\frac{C_r - N_r}{T_r}\right)} \\ T_n &= \frac{C_r \cdot T_r}{C_r - N_r} \quad \text{Equation (3)}\end{aligned}$$

where

$T_n$  = the new period of task,

To prevent an unbounded rise in the execution period of the task,  $T_n$ ,  $N_r$  must be limited to a maximum loan amount where  $N_r \ll C_r$ . If it is assumed that at its maximum value,  $T_n$ , is a multiple of  $T_r$ , then:

$$m \cdot T_r = \frac{C_r \cdot T_r}{C_r - N_m}$$

$$m = \frac{C_r \cdot T_r}{(C_r - N_m) \cdot T_r}$$

$$m = \frac{Cr}{Cr - Nm}$$

$$Cr - Nm = \frac{Cr}{m}$$

$$Nm = Cr - \frac{Cr}{m}$$

$$Nm = Cr \left( 1 - \frac{1}{m} \right) \quad \text{Equation (4)}$$

5 where

$m$  = multiple of the period of task<sub>r</sub>, and

$Nm$  = maximum execution time,  $Nr$ , loanable from task<sub>r</sub>.

FIG. 1 illustrates a real-time computing system 100 in accordance with the present invention. As shown in FIG. 1, the real-time computing system 100 includes certain  
10 standard hardware components, such as a processor 110 and a data storage device 120, such as a read-only memory and/or a random access memory (RAM).

The data storage device 120 includes a capacity loaning process 200, discussed further below in conjunction with FIG. 2. Generally, the capacity loaning process 200 implements capacity loaning between two tasks in accordance with the present invention.  
15 As shown in FIG. 1, the data storage device 120 also includes an operating system 150 that, among other things, manages the pairing of the two tasks, task<sub>U</sub> and task<sub>R</sub>, for capacity loaning in accordance with the present invention.

FIG. 2 illustrates a state diagram of the capacity loaning process 200. As shown in FIG. 2, each task is modeled here as a sporadic server, with modifications due to  
20 the capacity loan algorithm shown in bold type. For a discussion of the UML notations used herein, see, for example, UML Notation Guide, version 1.1 (Sept. 1, 1997), downloadable from <http://www.omg.org/docs/ad/97-08-05.pdf>, and incorporated by reference herein.

The capacity loaning process 200 utilizes two assumptions. First, the amount  
25 of borrowed time required by the urgent task will be known at the beginning of the execution period of the task. This implies that each event must be marked with its worst-case execution time so a dynamic assessment of budget can be made by the task at runtime. Secondly, the loan is only viable during the period of the borrower.

At system initialization, both tasks must agree on the maximum capacity that  
30 will be made available to loan. Using Equation (4), the routine task would compute the

maximum amount of execution time to loan based on its original execution budget and a limiting factor to prevent an unbounded rise in its period. This loanable amount would then be registered with the urgent task, which would then compute the maximum amount of execution time it may borrow with Equation (2) (accounting for the differences in the loan amount due to the differing period of the two tasks).

At the beginning of each execution period, the urgent task computes the amount of execution time over budget during state 220 and sends a message 240 to the routine task containing this amount (which may be 0, in the event that it is not over budget). The urgent task then executes these events normally during state 220.

Because the loan is only viable during the period of the borrower, the routine task maintains two state variables. The amount of capacity requested is the amount of execution time that has yet to be acted on. That is, this borrowed time has not yet degraded the performance of the loaner. The amount of capacity on loan is borrowed time that is currently acting to degrade the loaner. The action of the routine task in response to the message 240 from the urgent task depends on the task's current state. While waiting during state 250, the overbudget amount from the urgent task is applied to capacity requested. Otherwise, the overbudget amount is added to capacity on loan, where it is then used with Equation (3) to set a new period for the task. The priority of the task must also be lowered to match the new task period. At the beginning of the execution period of the routine task during state 260, capacity on loan is set from capacity requested, and Equation (3) is used to compute the task's period. Whenever the period of the task is changed, the corresponding priority must be changed accordingly, to satisfy the Rate Monotonic scheduling algorithm.

It is to be understood that the embodiments and variations shown and described herein are merely illustrative of the principles of this invention and that various modifications may be implemented by those skilled in the art without departing from the scope and spirit of the invention.

## CLAIMS:

1. A method for sharing execution capacity among tasks executing in a real-time computing system 100 having a performance specification in accordance with Rate Monotonic Analysis (RMA), comprising the steps of:

pairing a higher priority task, task<sub>U</sub>, with a lower priority task, task<sub>R</sub>;

- 5 reallocating execution time from the lower priority task, task<sub>R</sub>, to the higher priority task, task<sub>U</sub>, during an overload condition; and  
increasing the period of the lower priority task, task<sub>R</sub>, to compensate for said reallocated execution time.

- 10 2. The method of claim 1, wherein an amount of said execution time available to loan from said lower priority task, task<sub>R</sub>, to said higher priority task, task<sub>U</sub>, is obtained as follows:

$$N_u = \frac{N_r \cdot T_u}{T_r}$$

where,

- 15  $N_r$  = amount of execution time to borrow from task<sub>r</sub>, where  $N_r < C_r$ ,  
 $T_r$  = period of task<sub>r</sub>, and  
 $T_u$  = period of task<sub>U</sub>.

- 20 3. The method of claim 1, wherein said increased period of the lower priority task, task<sub>r</sub>, is obtained as follows:

$$T_n = \frac{C_r \cdot T_r}{C_r - N_r}$$

where

- $C_r$  = worst-case task execution time of task<sub>r</sub>,  
 $T_r$  = period of task<sub>r</sub>, and  
25  $N_r$  = amount of execution time to borrow from task<sub>r</sub>, where  $N_r < C_r$ .

4. The method of claim 1, further comprising the step of limiting an amount of execution time,  $N_r$ , to borrow from said lower priority task,  $task_r$ , to a maximum loan amount where  $N_r \ll C_r$ , where

$C_r$  = worst-case task execution time of  $task_r$ , and

5  $N_r$  = amount of execution time to borrow from  $task_r$ .

5. The method of claim 4, wherein a maximum execution time,  $N_m$ , that may be borrowed from said lower priority task,  $task_r$ , is obtained as follows:

$$N_m = C_r \left( 1 - \frac{1}{m} \right)$$

10 where  $m$  is the multiple of the period of said lower priority task,  $task_r$ .

6. A method for allocating resources among tasks executing in a real-time computing system 100 having a performance specification in accordance with Rate Monotonic Analysis (RMA), comprising the steps of:

15 pairing a higher priority task,  $task_U$ , with a lower priority task,  $task_R$ ;  
providing a first resource allocation to said lower priority task,  $task_R$ , during a normal operating condition; and  
reallocating a portion of said first resource allocation from said lower priority task,  $task_R$ , to said higher priority task,  $task_U$ , when said higher priority task,  $task_U$ , is operable.

20

7. A method for sharing execution capacity among tasks executing in a real-time computing system 100 having a performance specification in accordance with Rate Monotonic Analysis (RMA), comprising the steps of:

pairing a higher priority task,  $task_U$ , with a lower priority task,  $task_r$ ;  
25 reallocating execution time from the lower priority task,  $task_R$ , to the higher priority task,  $task_U$ , during an overload condition; and  
increasing the utilization of said higher priority task,  $task_U$ ; and  
decreasing the utilization of said lower priority task,  $task_R$ , in a proportional manner to maintain a constant utilization,  $U$ .

30

8. The method of claim 7, wherein said utilizations of said tasks are varied as follows:

$$\frac{C_u}{T_u} + \frac{C_r}{T_r} = U$$

where,

$C_u$  = worst-case task execution time of task<sub>u</sub>,

$T_u$  = period of task<sub>u</sub>,

5  $C_r$  = worst-case task execution time of task<sub>r</sub>,

$T_r$  = period of task<sub>r</sub>, and

$U$  = utilization for both tasks.

9. A real-time computing system 100 having a performance specification in  
 10 accordance with Rate Monotonic Analysis (RMA), comprising:  
 a memory 120 for storing computer readable code; and  
 a processor 110 operatively coupled to said memory 120, said processor 110 configured to:  
 pair a higher priority task, task<sub>u</sub>, with a lower priority task, task<sub>r</sub>;  
 reallocate execution time from the lower priority task, task<sub>r</sub>, to the higher priority task, task<sub>u</sub>,  
 15 during an overload condition; and  
 increase the period of the lower priority task, task<sub>r</sub>, to compensate for said reallocated  
 execution time.

10. A real-time computing system 100 having a performance specification in  
 20 accordance with Rate Monotonic Analysis (RMA), comprising:  
 a memory 120 for storing computer readable code; and  
 a processor 110 operatively coupled to said memory 120, said processor 110 configured to:  
 pair a higher priority task, task<sub>u</sub>, with a lower priority task, task<sub>r</sub>;  
 provide a first resource allocation to said lower priority task, task<sub>r</sub>, during a normal operating  
 25 condition; and  
 reallocate a portion of said first resource allocation from said lower priority task, task<sub>r</sub>, to  
 said higher priority task, task<sub>u</sub>, when said higher priority task, task<sub>u</sub>, is operable.

11. A real-time computing system 100 having a performance specification in  
 30 accordance with Rate Monotonic Analysis (RMA), comprising:  
 a memory 120 for storing computer readable code; and  
 a processor 110 operatively coupled to said memory 120, said processor 110 configured to:  
 pair a higher priority task, task<sub>u</sub>, with a lower priority task, task<sub>r</sub>;

reallocate execution time from the lower priority task,  $\text{task}_R$ , to the higher priority task,  $\text{task}_U$ , during an overload condition; and

increase the utilization of said higher priority task,  $\text{task}_U$ ; and

decrease the utilization of said lower priority task,  $\text{task}_R$ , in a proportional manner to maintain

5 a constant utilization,  $U$ .

REAL-TIME COMPUTING SYSTEM-100

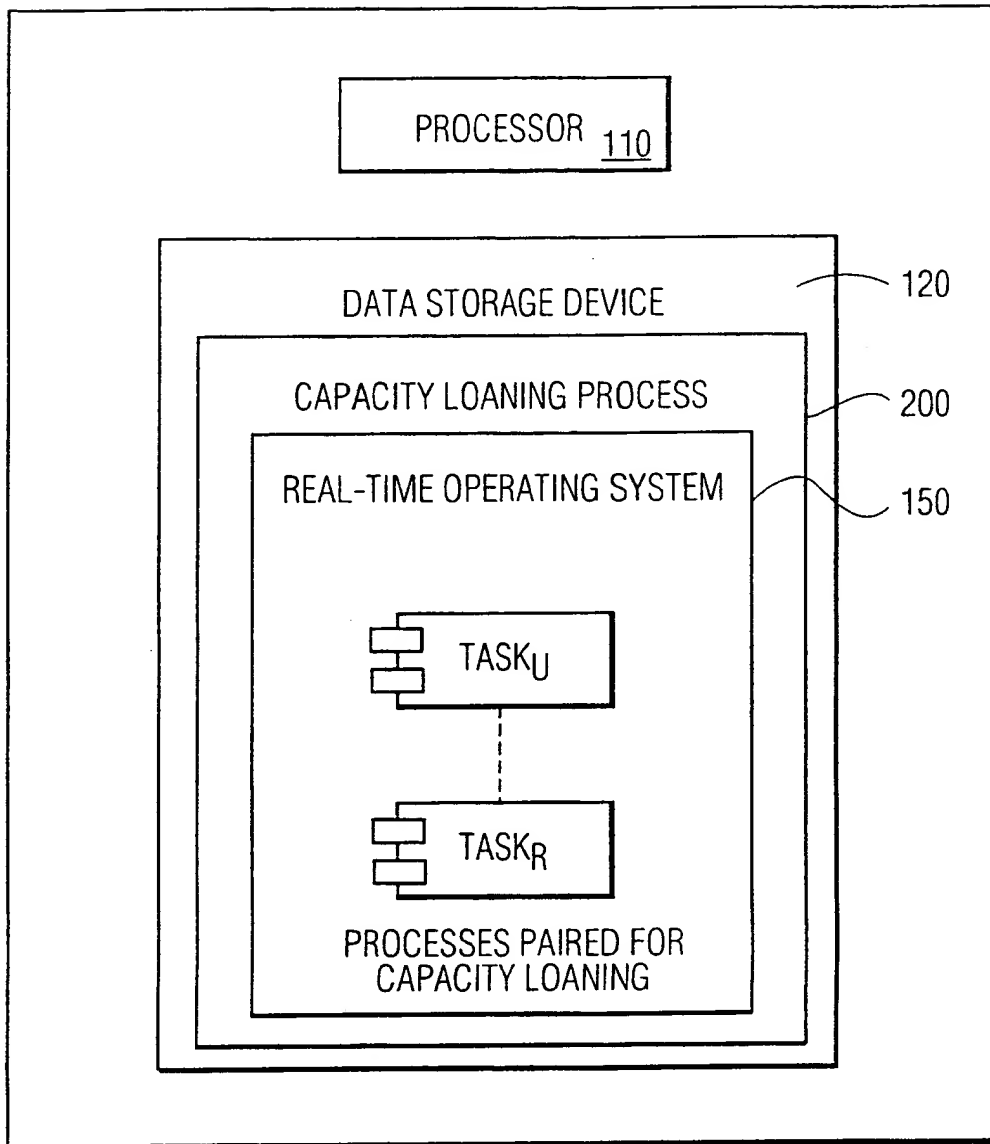


FIG. 1



## CAPACITY LOANING PROCESS-200

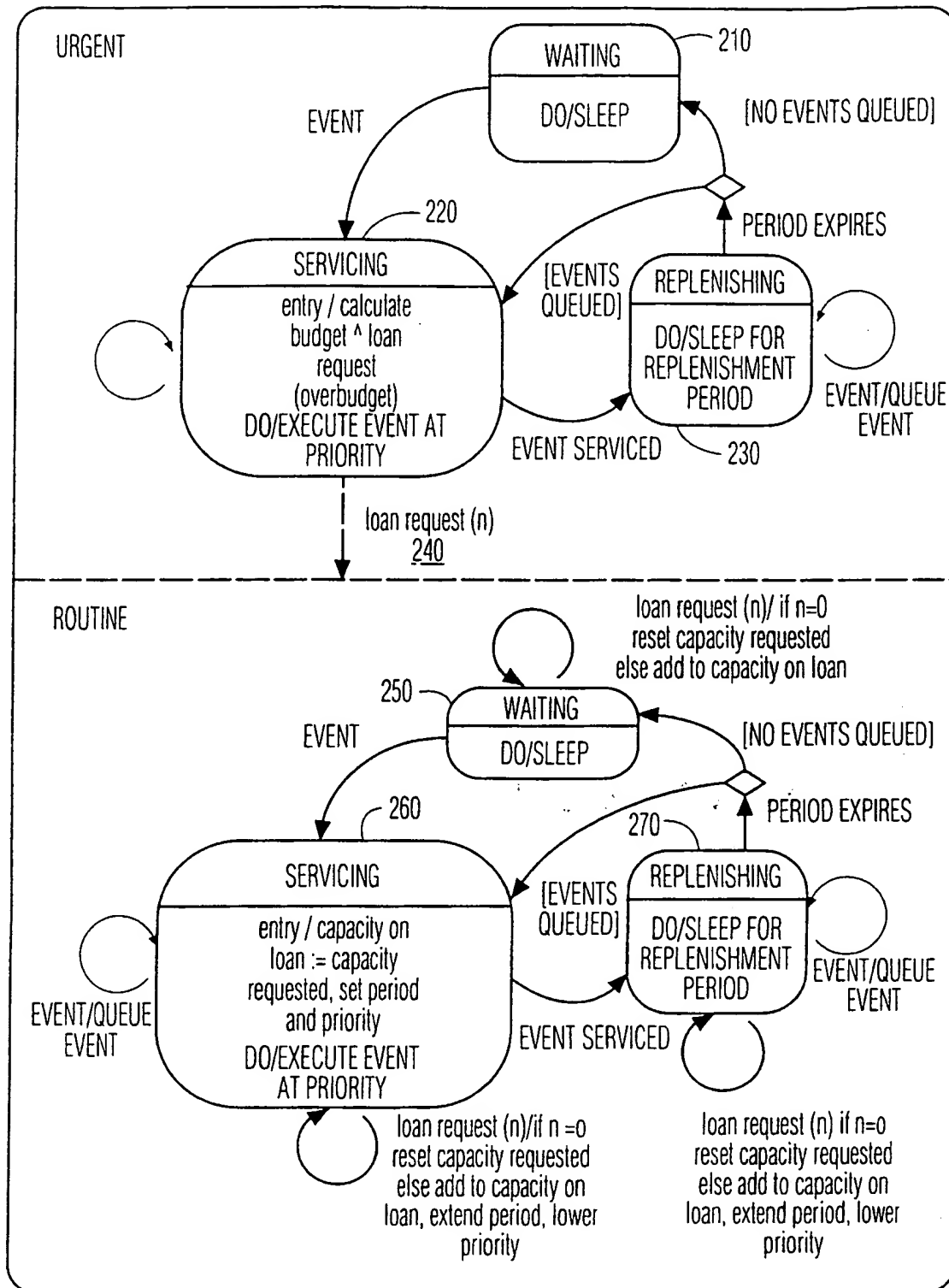


FIG. 2

**THIS PAGE BLANK (USPTO)**

(19) World Intellectual Property Organization  
International Bureau

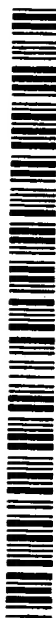


(43) International Publication Date  
19 October 2000 (19.10.2000)

PCT

(10) International Publication Number  
**WO 00/62157 A3**

- (51) International Patent Classification<sup>7</sup>: G06F 9/50, 9/48 (74) Agent: GRAVENDEEL, Cornelis; Internationaal Octrooibureau B.V., Prof. Holstlaan 6, NL-5656 AA Eindhoven (NL).
- (21) International Application Number: PCT/EP00/03204
- (22) International Filing Date: 11 April 2000 (11.04.2000) (81) Designated State (*national*): JP.
- (25) Filing Language: English (84) Designated States (*regional*): European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).
- (26) Publication Language: English
- (30) Priority Data:  
60/129,301 14 April 1999 (14.04.1999) US  
09/481,771 11 January 2000 (11.01.2000) US
- (71) Applicant: KONINKLIJKE PHILIPS ELECTRONICS N.V. [NL/NL]; Groenewoudseweg 1, NL-5621 BA Eindhoven (NL).
- (72) Inventor: ISHAM, Karl, M.; Prof. Holstlaan 6, NL-5656 AA Eindhoven (NL).
- Published:  
— With international search report.  
— Before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments.
- (88) Date of publication of the international search report:  
8 February 2001
- For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.



WO 00/62157 A3

(54) Title: METHOD FOR DYNAMIC LOANING IN RATE MONOTONIC REAL-TIME SYSTEMS

(57) Abstract: A method and apparatus are disclosed for sharing execution capacity among tasks executing in a real-time computing system. The present invention extends RMA techniques for characterizing system timing behavior and designing real-time systems. A high priority task having hard deadlines is paired with a lower priority task having soft deadlines. During an overload condition, the higher priority task can dynamically borrow execution time from the execution capacity of the lower priority task without affecting the schedulability of the rest of the system. The higher priority task is bolstered in a proportion to the capacity borrowed from the lower priority task, so that the combined utilization of the two tasks remains constant. The period of the degraded task is increased to compensate for the execution time that was loaned to the higher priority task. In addition, the priority of the lower priority task is modified to match the new period.

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> IPC 7 G06F9/50 G06F9/48				
According to International Patent Classification (IPC) or to both national classification and IPC				
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) IPC 7 G06F				
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched				
Electronic data base consulted during the international search (name of data base and, where practical, search terms used) IBM-TDB, EPO-Internal, PAJ				
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>				
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.		
X	LUI SHA ET AL: "MODE CHANGE PROTOCOLS FOR PRIORITY-DRIVEN PREEMPTIVE SCHEDULING" REAL TIME SYSTEMS, NL, KLUWER ACADEMIC PUBLISHERS, DORDRECHT, vol. 1, no. 3, 1 December 1989 (1989-12-01), pages 243-264, XP000281989 ISSN: 0922-6443 page 245, paragraph 2.1 -page 246 page 258, line 32 -page 261, line 22 --- -/--	1,4,6-11		
<input checked="" type="checkbox"/> Further documents are listed in the continuation of box C. <input type="checkbox"/> Patent family members are listed in annex.				
* Special categories of cited documents:				
<table style="width: 100%; border: none;"> <tr> <td style="width: 50%; vertical-align: top; border: none;">           *A* document defining the general state of the art which is not considered to be of particular relevance            *E* earlier document but published on or after the international filing date            *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)            *O* document referring to an oral disclosure, use, exhibition or other means            *P* document published prior to the international filing date but later than the priority date claimed         </td> <td style="width: 50%; vertical-align: top; border: none;">           *T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention            *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone            *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.            *Z* document member of the same patent family         </td> </tr> </table>			*A* document defining the general state of the art which is not considered to be of particular relevance *E* earlier document but published on or after the international filing date *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) *O* document referring to an oral disclosure, use, exhibition or other means *P* document published prior to the international filing date but later than the priority date claimed	*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art. *Z* document member of the same patent family
*A* document defining the general state of the art which is not considered to be of particular relevance *E* earlier document but published on or after the international filing date *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) *O* document referring to an oral disclosure, use, exhibition or other means *P* document published prior to the international filing date but later than the priority date claimed	*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art. *Z* document member of the same patent family			
Date of the actual completion of the international search  <div style="text-align: center; font-weight: bold;">23 November 2000</div>	Date of mailing of the international search report  <div style="text-align: center; font-weight: bold;">30/11/2000</div>			
Name and mailing address of the ISA European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Tx. 31 651 epo nl, Fax: (+31-70) 340-3016	Authorized officer  <div style="text-align: center; font-weight: bold;">Michel, T</div>			

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>SILVERTHORN L: "RATE-MONOTONIC SCHEDULING ENSURES TASKS MEET DEADLINES" EDN ELECTRICAL DESIGN NEWS,US,CAHNERS PUBLISHING CO. NEWTON, MASSACHUSETTS, vol. 34, no. 22, 26 October 1989 (1989-10-26), pages 191-198,200, XP000070837 ISSN: 0012-7515 page 192, left-hand column, line 35 -right-hand column, line 21 page 198.2, left-hand column, line 5 -right-hand column, line 5</p>	1,4,6-11
A	<p>SHA L ET AL: "GENERALIZED RATE-MONOTONIC SCHEDULING THEORY: A FRAMEWORK FOR DEVELOPING REAL-TIME SYSTEMS" PROCEEDINGS OF THE IEEE,US,IEEE. NEW YORK, vol. 82, no. 1, 1994, pages 68-82, XP000435882 ISSN: 0018-9219 page 69, paragraph A -page 70 page 71, paragraph C -page 72</p>	1,4,6-11
A	<p>STREICH H: "TASKPAIR-SCHEDULING: AN APPROACH FOR DYNAMIC REAL-TIME SYSTEMS" INTERNATIONAL JOURNAL OF MINI AND MICROCOMPUTERS,US,ACTA PRESS. ANAHEIM, CALIFORNIA, vol. 17, no. 2, 1995, pages 77-83, XP000518020 ISSN: 0702-0481 abstract page 78, left-hand column, line 22 - last line</p>	1,6,7, 9-11
A	<p>SHIH W K ET AL: "MODIFIED RATE-MONOTONIC ALGORITHM FOR SCHEDULING PERIODIC JOBS WITH DEFERRED DEADLINES" IEEE TRANSACTIONS ON SOFTWARE ENGINEERING,US,IEEE INC. NEW YORK, vol. 19, no. 12, 1 December 1993 (1993-12-01), pages 1171-1179, XP000418805 ISSN: 0098-5589 the whole document</p>	1,6,7, 9-11

**THIS PAGE BLANK (USPTO)**

**THIS PAGE BLANK (USPTO)**

**THIS PAGE BLANK (USPTO)**

**(PTO)**

**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

**BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☒ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**

**THIS PAGE BLANK (USPTO)**